

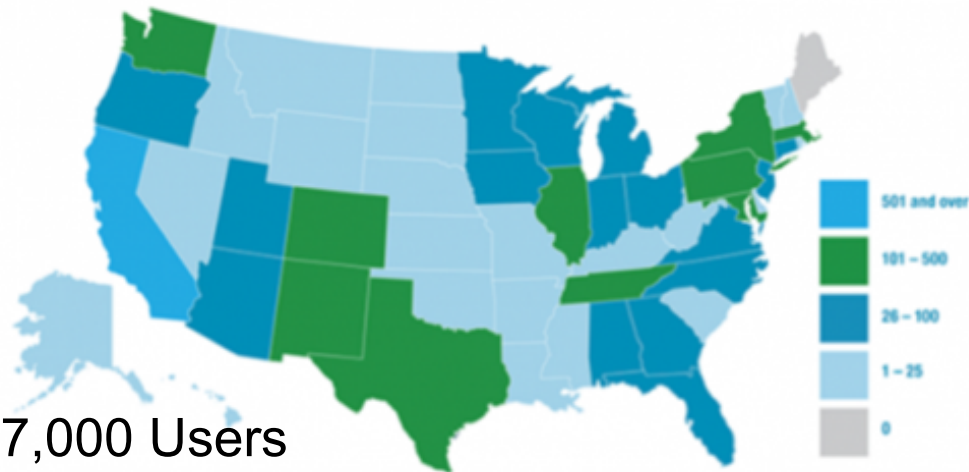
Upcoming computing facilities for science



Snowmass Computational Frontier
Workshop

Nicholas J Wright
NERSC Chief Architect
11 August, 2020

NERSC is the mission High Performance Computing facility for the DOE SC

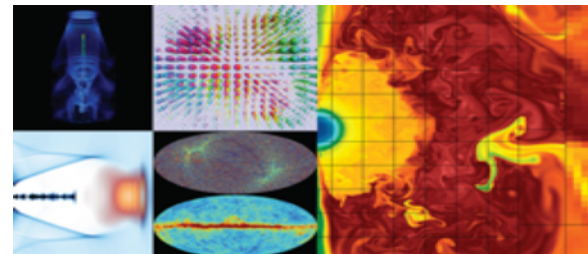
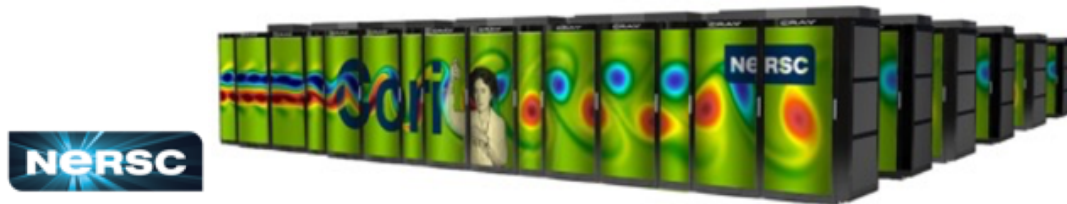


7,000 Users

800 Projects

700 Codes

2000 NERSC citations per year



Simulations at scale



Data analysis support for
DOE's experimental and
observational facilities

Photo Credit: CAMERA



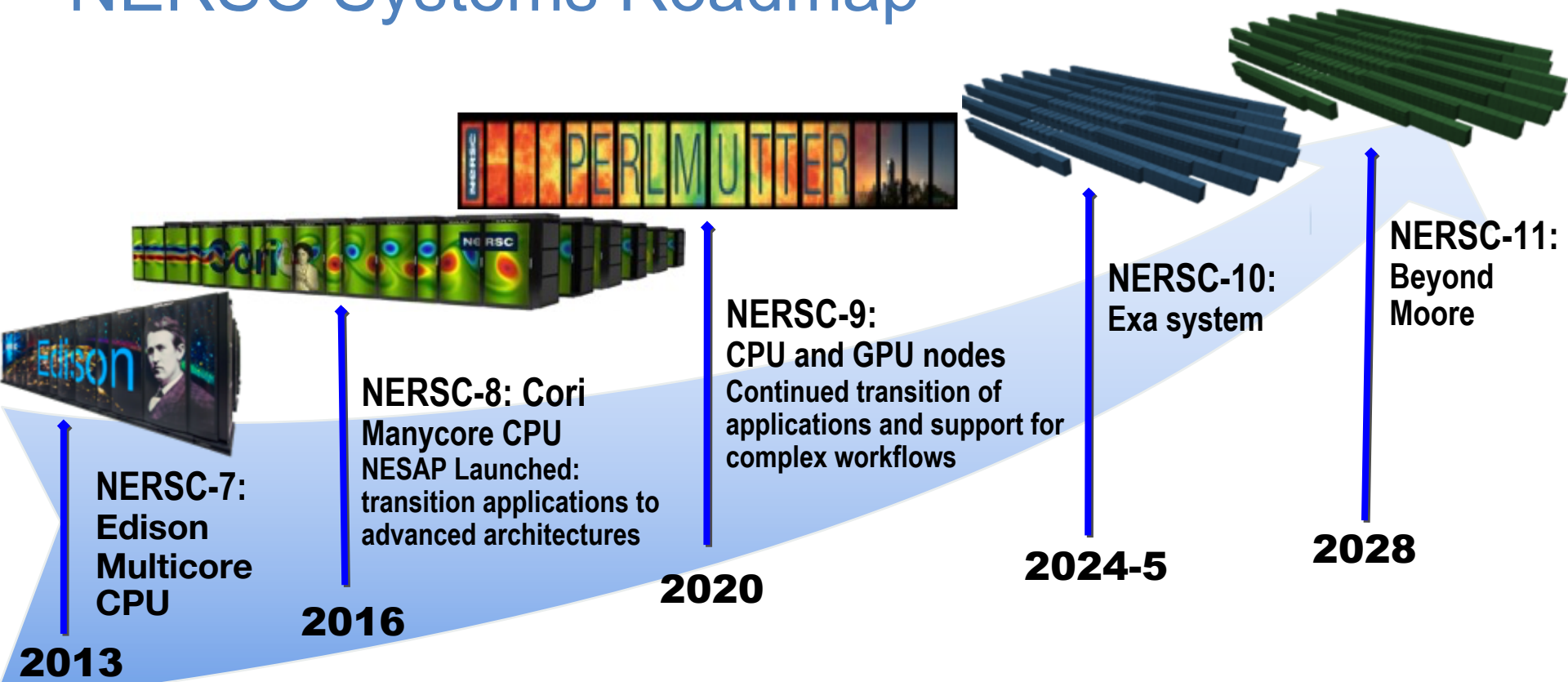
BERKELEY LAB
Bringing Science Solutions to the World



U.S. DEPARTMENT OF
ENERGY

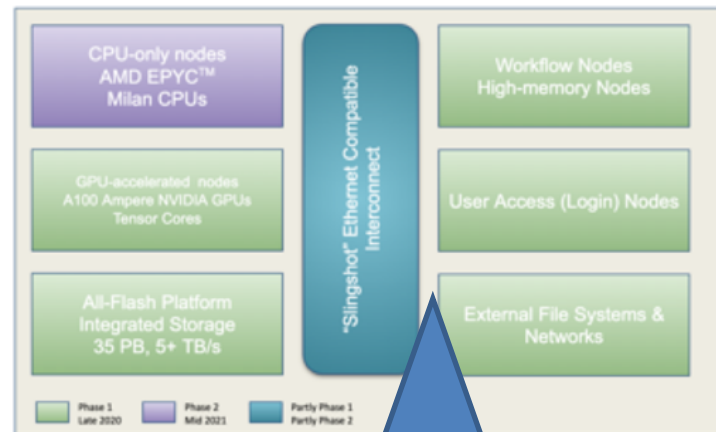
Office of
Science

NERSC Systems Roadmap



Perlmutter: A System Optimized for Science

- Cray Shasta System providing 3-4x capability of Cori
- GPU-accelerated and CPU-only for large scale simulation and data analysis from experimental facilities
- GPU nodes: 4 NVIDIA A100 “Ampere” GPUs each w/Tensor Cores & NVLink-3 and High-BW memory + 1 AMD “Milan” CPU
 - Over 6000 NVIDIA Volta-Next GPUs
 - Unified Virtual Memory support improves programmability
- Cray “Slingshot” - High-performance, scalable, low-latency Ethernet- compatible network
 - Capable of Terabit connections to/from the system
- Single-tier All-Flash Lustre based HPC file system



Phased
delivery
1st phase: End
CY2020
2nd phase:
Spring CY2021

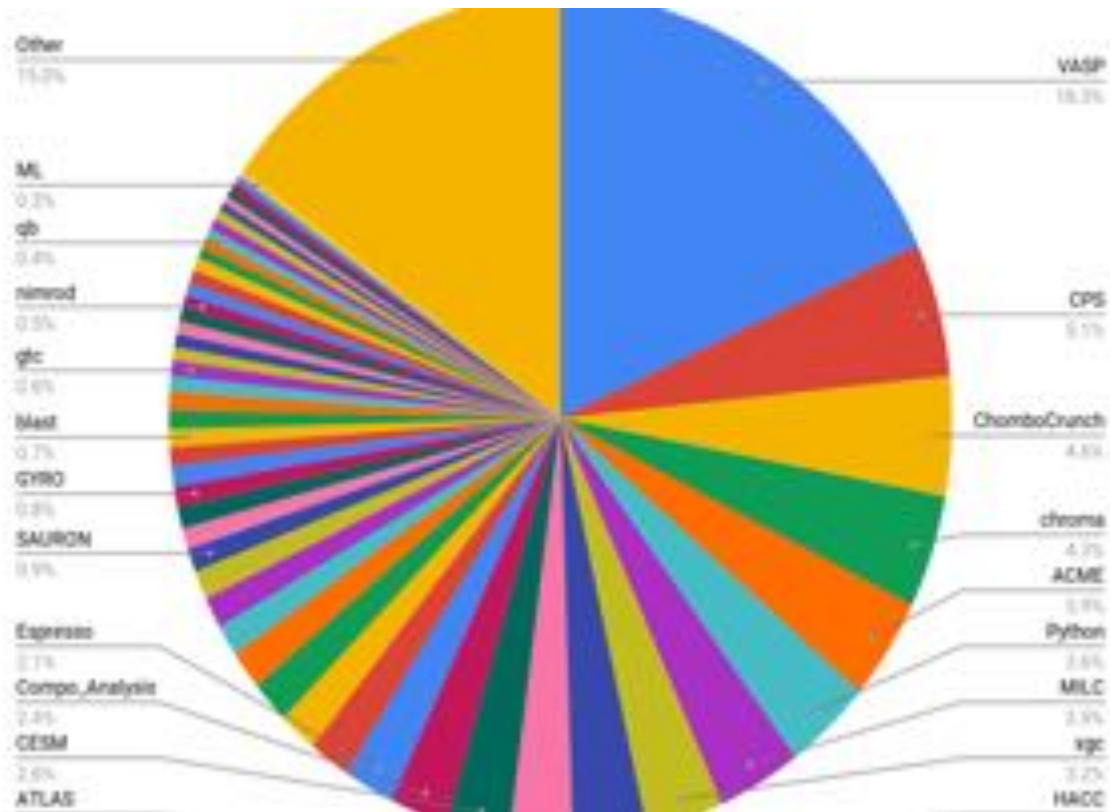


BERKELEY LAB
Bringing Science Solutions to the World



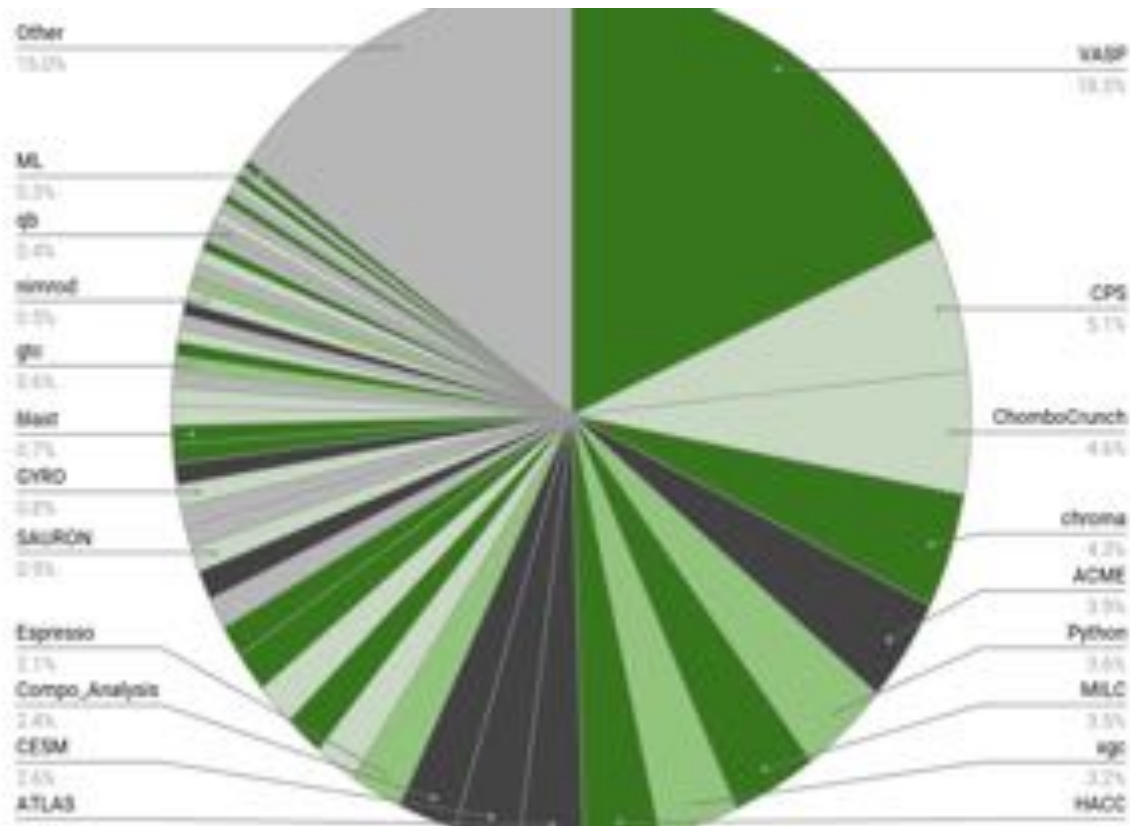
U.S. DEPARTMENT OF
ENERGY | Office of
Science

NERSC System Utilization (Aug'17 - Jul'18)



- 3 codes > 25% of the workload
- 10 codes > 50% of the workload
- 35 codes > 75% of the workload
- Over 600 codes comprise the remaining 25% of the workload.

GPU Readiness Among NERSC Codes (Aug'17 - Jul'18)



| GPU Status & Description | Fraction |
|---|----------|
| Enabled: Most features are ported and performant | 37% |
| Kernels: Ports of some kernels have been documented. | 10% |
| Proxy: Kernels in related codes have been ported | 20% |
| Unlikely: A GPU port would require major effort. | 13% |
| Unknown: GPU readiness cannot be assessed at this time. | 20% |



A number of applications in NERSC workload are GPU enabled already.



BERKELEY LAB
Bringing Science Solutions to the World



U.S. DEPARTMENT OF
ENERGY
Office of Science

Pre-Exascale Systems

2013

2016

2018

2020

Exascale Systems

2021-2023



NVIDIA



Argonne
Intel/Cray



LBNL
Cray/NVIDIA/AMD



ORNL
Cray/AMD



LANL/SNL
TBD



LLNL
Cray/AMD



Compute Node

2 Intel Xeon scalable "Sapphire Rapids" processors; 6 Xe arch-based GPUs; Unified Memory Architecture; 8 fabric endpoints

CPU-GPU Interconnect

CPU-GPU: PCIe; GPU-GPU: Xe Link

Sustained Performance

≥ 1 Exaflop DP

Platform

Cray Shasta

Software Stack

Cray Shasta software stack + Intel enhancements + data and learning

System Interconnect

System Interconnect Cray Slingshot; Dragonfly topology with adaptive routing

High-Performance Storage

≥ 230 PB, ≥ 25 TB/s (DAOS)

Aggregate System Memory

> 10 PB

GPU Architecture

Xe arch-based "Ponte Vecchio" GPU; Tile-based chiplets, HBM stack, Foveros 3D integration, 7nm

Network Switch

25.6 Tb/s per switch, from 64–200 Gbs ports (25 GB/s per direction)

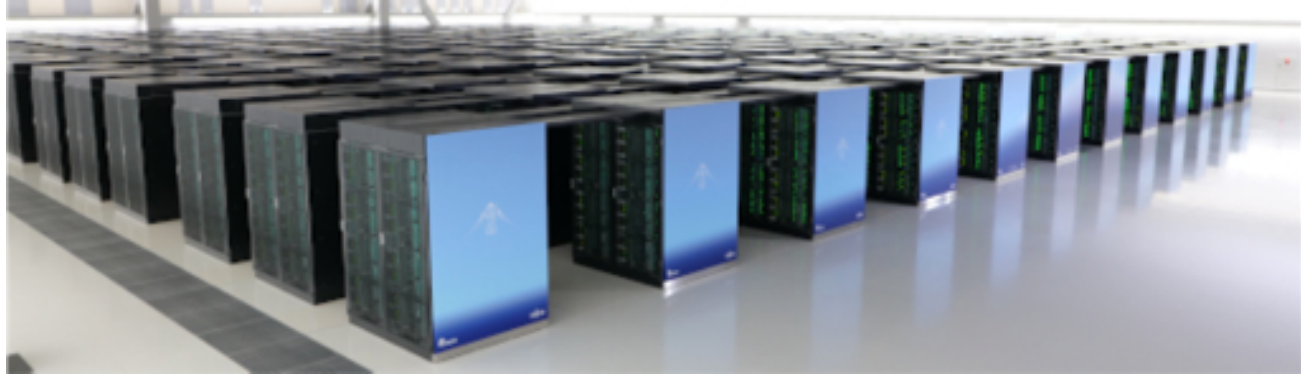
Programming Models

Intel oneAPI, MPI, OpenMP, C/C++, Fortran, SYCL/DPC++

| | |
|----------------------|--|
| Peak Performance | >1.5 EF |
| Footprint | > 100 cabinets |
| Node | 1 HPC and AI Optimized AMD EPYC CPU 4 Purpose Built AMD Radeon Instinct GPU |
| CPU-GPU Interconnect | AMD Infinity Fabric Coherent memory across the node |
| System Interconnect | Multiple Slingshot NICs providing 100 GB/s network bandwidth Slingshot dragonfly network which provides adaptive routing, congestion management and quality of service. |
| Storage | 2-4x performance and capacity of Summit's I/O subsystem. Frontier will have near node storage like Summit. |



Fugaku



| | |
|---------------------|--|
| Peak Performance | 488 PF |
| Footprint | 158,976 nodes, 414 racks |
| Node | 48 core ARM – 3 TF 32 GiB HBM 1,024 GB/s |
| System Interconnect | 5D TOFU |
| Storage | Every 16 Nodes local SSD 1.6 TB 150 PB Lustre |
| Power | 29 MW |



Further out.....

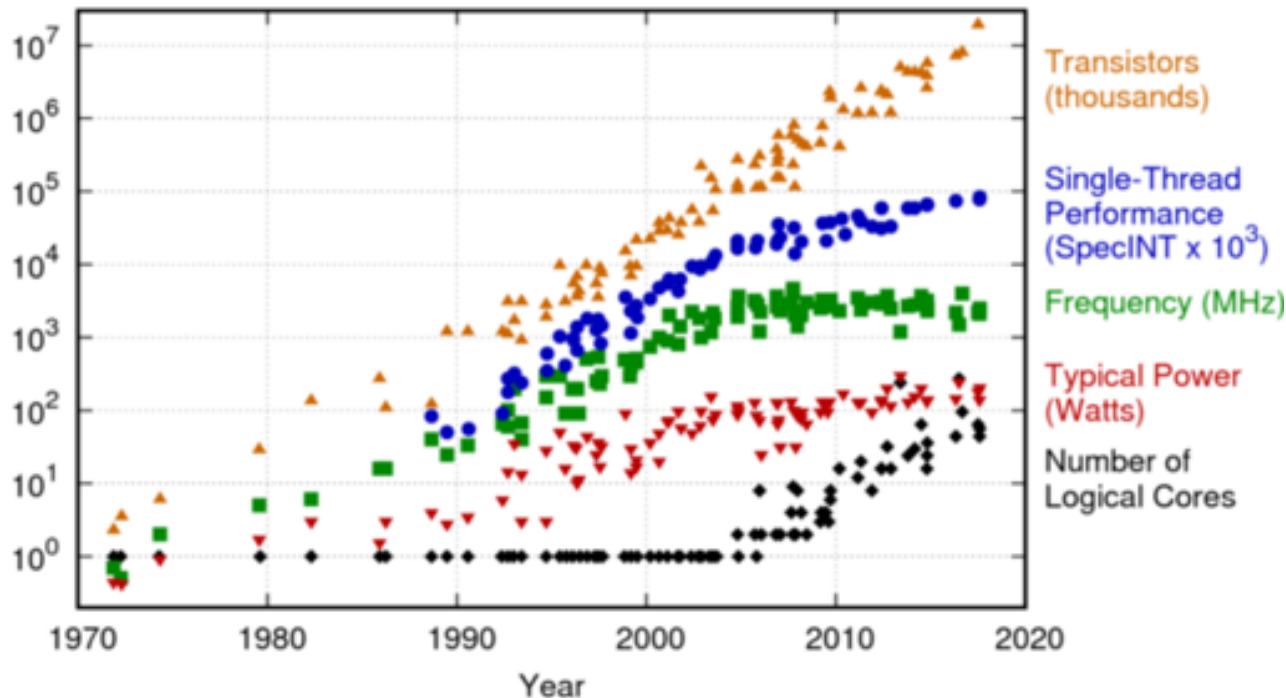
- Seems reasonable to assume all 3 DOE ASCR facilities will upgrade in the mid 20's
 - ~10's EF
- NSF
 - Frontera @ TACC ~ \$60M ~40 PF
 - 8,000 dual socket Xenon “Cascade Lake” 28 cores / socket, HDR-IB, 90 4-way NVIDIA GPU Nodes (Volta)
 - Phase II – 10x phase I (~400 PF) due 2023-24

What will these machines look like ?

Technology Scaling Trends

42 Years of Microprocessor Trend Data

Performance



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

> serial performance

> Performance per socket

- More cores / CUDA cores
- Multichip packages to keep costs down – NUMA

More heterogeneity

- CPU's and GPU's are here to stay
- AI accelerators ?
- Compute in network / storage

Innovations like domain-specific hardware, enhanced security, open instruction sets, and agile chip development will lead the way.

BY JOHN L. HENNESSY AND DAVID A. PATTERSON

A New Golden Age for Computer Architecture

Extreme Heterogeneity 2018

PRODUCTIVE COMPUTATIONAL SCIENCE
IN THE ERA OF EXTREME HETEROGENEITY



End of Moore's Law ?

EE|Times

HOME NEWS ▾ PERSPECTIVES DESIGNLINES ▾ VIDEOS RADIO EDUCATION ▾

DESIGNLINES | SOC DESIGNLINE

TSMC Aims to Build World's First 3-nm Fab

TSMC will build the world's first 3-nm fab the company does the bulk of its

ANANDTECH

PC COMPONENTS ▾ SMARTPHONES & TABLETS ▾ SYSTEMS ▾ ENTERPRISE

TRENDING TOPICS CPU INTEL AMD MOBILE SMARTPHONES STORAGE RYZEN SSDS GPU

Home > Semiconductors

Samsung Announces 3nm GAA MBCFET PDK, Version 0.1

by Ian Cutress on May 14, 2018 8:00 PM EST

Posted in: Semiconductors, Samsung, 3nm, GAAFET, MBCFET



Planar FET



FinFET

ANANDTECH

PC COMPONENTS ▾ SMARTPHONES & TABLETS ▾ SYSTEMS ▾ ENTERPRISE

TRENDING TOPICS CPU INTEL AMD MOBILE SMARTPHONES STORAGE RYZEN SSDS GPU

Home > CPUs

Intel Details Manufacturing through 2023: 7nm, 7+, 7++, with Next Gen Packaging

by Ian Cutress & Anton Shilov on May 8, 2018 4:35 PM EST

Posted in: CPUs, Intel, 7nm, 7nm+, 7nm++, EMB, 7nm+, POWERPC, Xeon, Xeon+, Xeon+, Zen, Zen+, SPARC



BERKELEY LAB
Bringing Science Solutions to the World

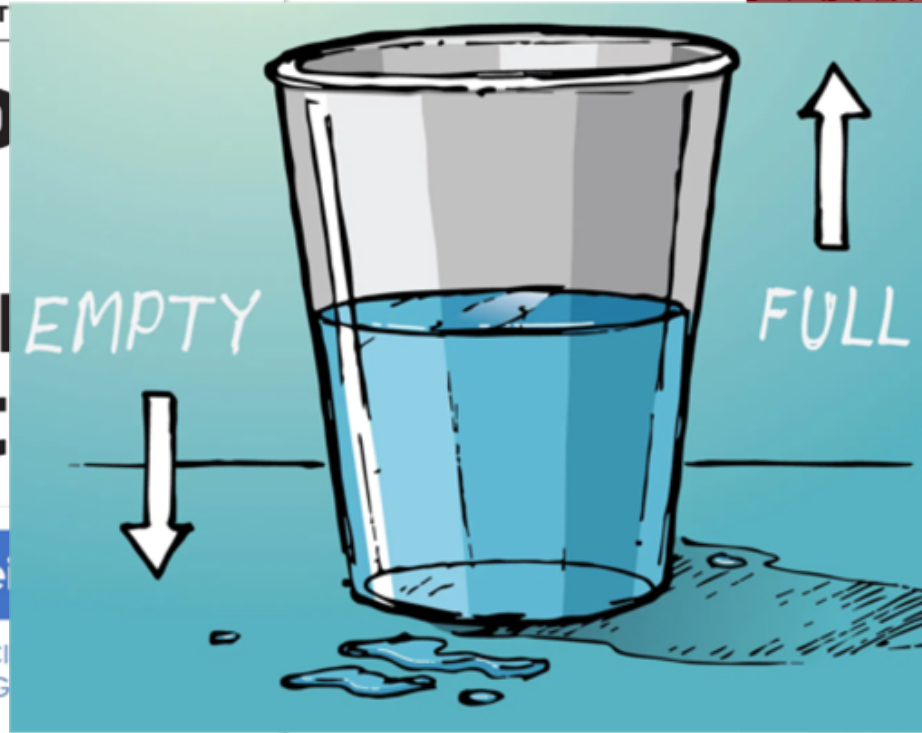


U.S. DEPARTMENT OF
ENERGY | Office of
Science

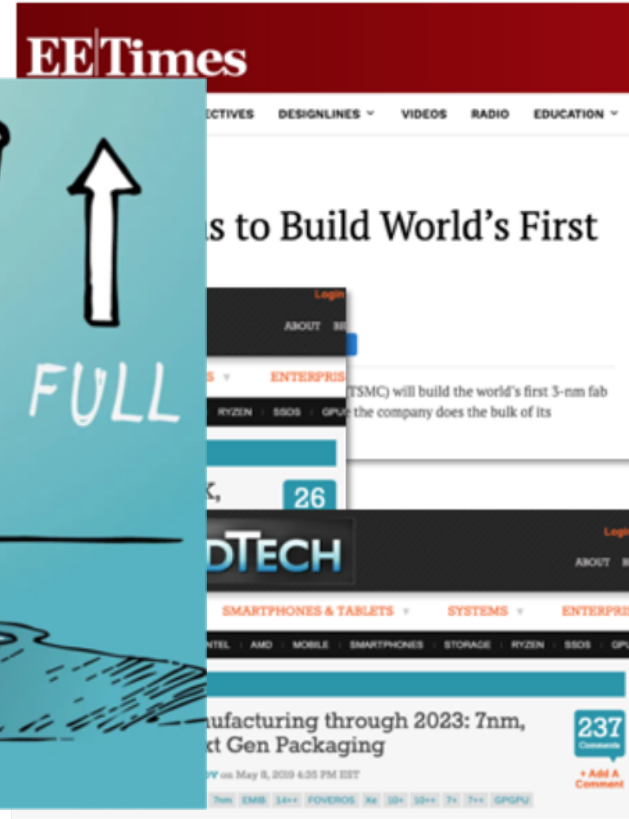
A New Go Age for Computer Architect

Extreme Heterogeneity

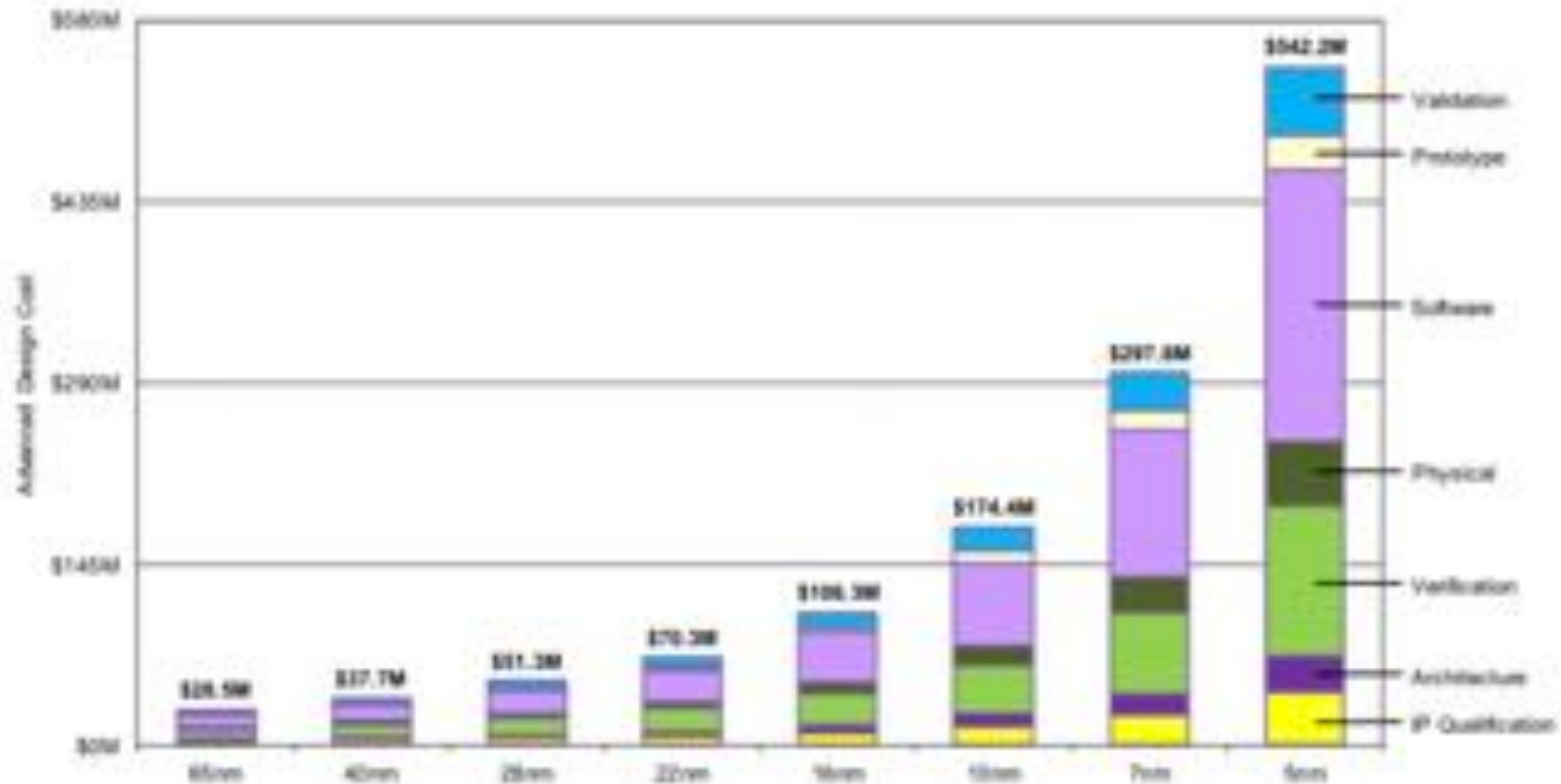
PRODUCTIVE COMPUTATIONAL SCIENCE IN THE ERA OF EXTREME HETEROGENEITY



End of Moore's Law ?



Can I just do it myself ? - Chip Design Costs



Summary

- Future HPC resources are likely to contain CPUs & GPUs
 - They continue to increase in performance each generation
 - Almost certainly on a per \$ basis
 - Maybe not as much per Watt ~1.5x per generation
 - Different sites may continue to have differing amounts of each - although more tightly-coupled solutions may affect this
- Lifetime of systems will increase as incentive to upgrade gets less
 - Expect more specialization

Questions ?



Will GPUs work for everybody?

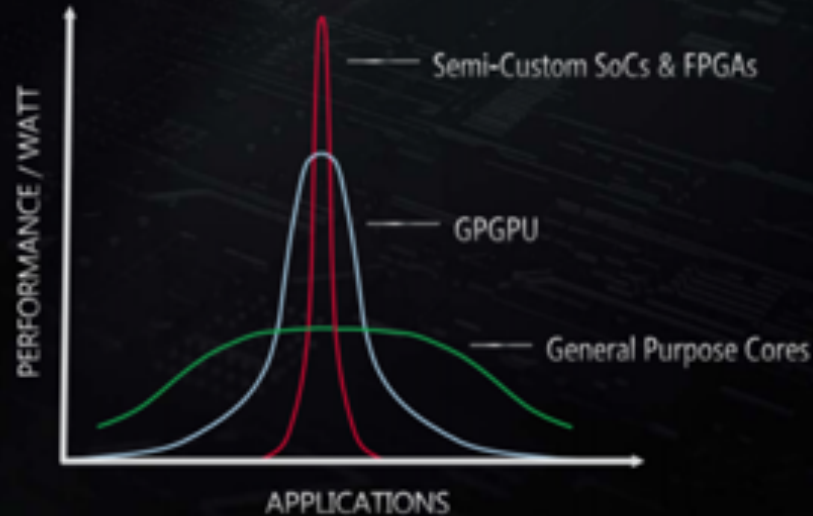
- Will 100% of the NERSC workload be able to utilize GPUs by 2024?
 - Yes, they just need to modify their code
 - No, their algorithm needs changing
 - No, their physics is fundamentally not amenable to data parallelism
 - No, they just don't have time or need too



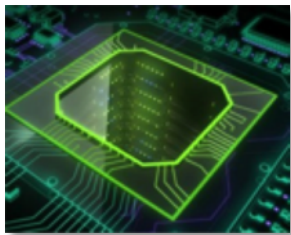
View from AMD - can we exploit this to benefit NERSC users ?



OPTIMIZING SYSTEM PERFORMANCE WITH HETEROGENEOUS COMPUTING



Specialization: End Game for Moore's Law



NVIDIA builds deep learning appliance with V100 Tesla's



RISC-V is an open hardware platform



Intel buys deep learning startup, Nervana

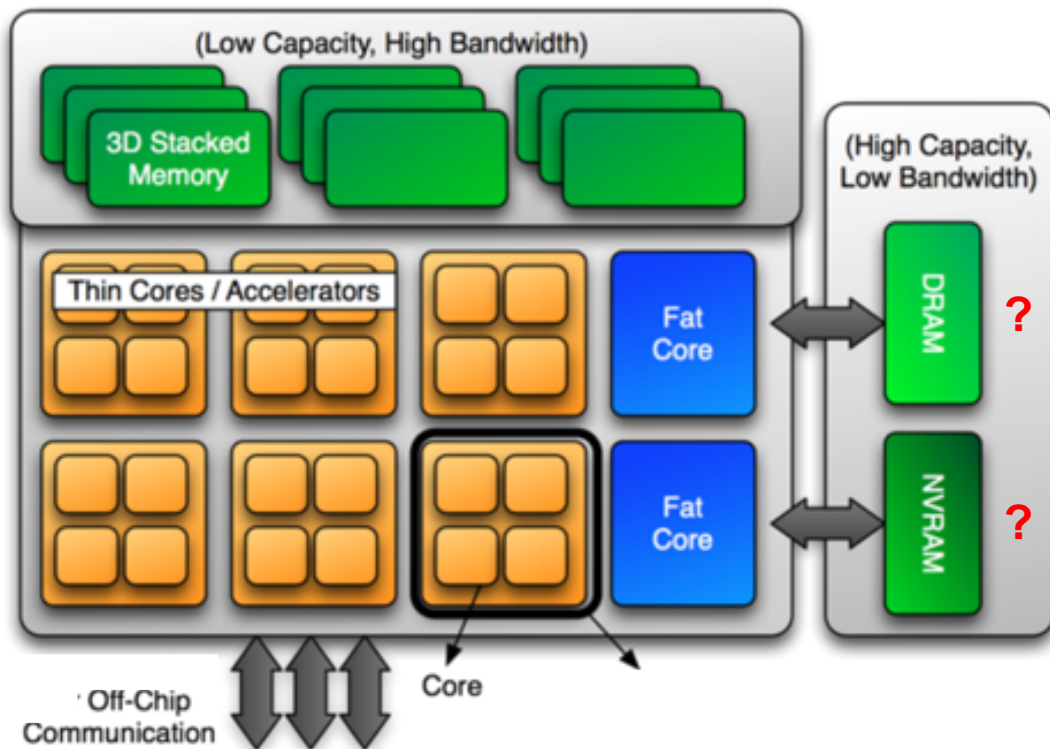


FPGAs offer configurable specialization



Google designs its own Tensor Processing Unit (TPU)

Potential 2024 Node



- Vendors converging to a mixture of energy-efficient Thin Cores/Accelerators and Fat Cores
- Potentially with DRAM/NVRAM
- (Hopefully) leads to less focus on data motion and more on identifying parallelism